

# Space Exploration via Proximity Search

Sariel Har-Peled\*   Nirman Kumar†   David M. Mount‡   Benjamin Raichel

December 4, 2014

## Abstract

We investigate what computational tasks can be performed on a point set in  $\mathbb{R}^d$ , if we are only given black-box access to it via nearest-neighbor search. This is a reasonable assumption if the underlying point set is either provided implicitly, or it is stored in a data structure that can answer such queries. In particular, we show the following:

- (A) One can compute an approximate bi-criteria  $k$ -center clustering of the point set, and more generally compute a greedy permutation of the point set.
- (B) One can decide if a query point is (approximately) inside the convex-hull of the point set.

We also investigate the problem of clustering the given point set, such that meaningful proximity queries can be carried out on the centers of the clusters, instead of the whole point set.

## 1. Introduction

Many problems in Computational Geometry involve sets of points in  $\mathbb{R}^d$ . Traditionally, such a point set is presented explicitly, say, as a list of coordinate vectors. There are, however, numerous applications in science and engineering where point sets are presented *implicitly*. This may arise for various reasons: (1) the point set (which might be infinite) is a physical structure that is represented in terms of a finite set of sensed measurements such as a point cloud, (2) the set is too large to be stored explicitly in memory, or (3) the set is procedurally generated from a highly compressed form. (A number of concrete examples are described below.)

Access to such an implicitly-represented point set  $P$  is performed through an *oracle* that is capable of answering queries of a particular type. We can think of this oracle as a black-box data structure, which is provided to us in lieu of an explicit representation. Various types of probes have been studied (such as finger probes, line probes, and X-ray probes [Ski89]). Most of these assume that  $P$  is connected (e.g., a convex polygon) and cannot be applied when dealing with arbitrary point sets. In this paper we

---

\*Department of Computer Science; University of Illinois; 201 N. Goodwin Avenue; Urbana, IL, 61801, USA; [sariel@illinois.edu](mailto:sariel@illinois.edu); <http://sarielhp.org/>. Work on this paper was partially supported by a NSF AF award CCF-0915984.

†Department of Computer Science; University of California; 2120B Harold Frank Hall; Santa Barbara, CA, 93106, USA; [nirman@cs.ucsb.edu](mailto:nirman@cs.ucsb.edu); Work on this paper was partially supported by a NSF AF award CCF-1217462 while the author was a student at UIUC, and by NSF grant CCF-1161495 and a grant from DARPA while the author has been a postdoc at UCSB.

‡Department of Computer Science; University of Maryland; College Park, MD, 20742, USA; [mount@cs.umd.edu](mailto:mount@cs.umd.edu); <http://www.cs.umd.edu/~mount/>. Work on this paper was partially supported by NSF award CCF-1117259 and ONR award N00014-08-1-1015.

consider a natural choice for probing general point sets based on computing nearest neighbors, which we call *proximity probes*.

More formally, we assume that the point set  $P$  is a (not necessarily finite) compact subset of  $\mathbb{R}^d$ . The point set  $P$  is accessible only through a nearest-neighbor data structure, which given a query point  $q$ , returns the closest point of  $P$  to  $q$ . Some of our results assume that the data structure returns an exact nearest neighbor (NN) and others assume that the data structure returns a  $(1 + \varepsilon)$ -approximate nearest-neighbor (ANN). (See [Section 2](#) for definitions.) In any probing scenario, it is necessary to begin with a general notion of the set’s spatial location. We assume that  $P$  is contained within a given compact subset  $\mathcal{D}$  of  $\mathbb{R}^d$ , called the *domain*.

The oracle is given as a black-box. Specifically, we do not allow deletions from or insertions into the data structure. We do not assume knowledge of the number of data points, nor do we make any continuity or smoothness assumptions. Indeed, most of our results apply to infinite point sets, including volumes or surfaces.

## Prior Work and Applications

Implicitly-represented point sets arise in various applications. One example is that of analyzing a geometric shape through probing. An example of this is Atomic Force Microscopy (AFM) [\[Wik14\]](#). This technology can reveal the undulations of a surface at the resolution of fractions of a nanometer. It relies on the principle that when an appropriately designed tip (the probe) is brought in the proximity of a surface to scan it, certain atomic forces minutely deflect the tip in the direction of the surface. Since the deflection of the tip is generally to the closest point on the surface, this mode of acquisition is an example of proximity probing. A sufficient number of such samples can be used to reconstruct the surface [\[BQG86\]](#).

The topic of shape analysis through probing has been well studied within the field of computational geometry. The most commonly assumed probe is a *finger probe*, which determines the first point of contact of a ray and the set. Cole and Yap [\[CY87\]](#) pioneered this area by analyzing the minimum number of finger probes needed to reconstruct a convex polygon. Since then, various alternative probing methods have been considered. For good surveys of this area, see Skiena [\[Ski89, Ski97\]](#).

More recently, Boissonnat *et al.* [\[BGO07\]](#) presented an algorithm for learning a smooth unknown surface  $S$  bounding an object  $\mathcal{O}$  in  $\mathbb{R}^3$  through the use of finger probes. Under some reasonable assumptions, their algorithm computes a triangulated surface  $\hat{S}$  that approximates  $S$  to a given level of accuracy. In contrast to our work, which applies to general point sets, all of these earlier results assume that the set in question is a connected shape or surface.

Implicitly-represented point sets also arise in geometric modeling. Complex geometric sets are often generated from much smaller representations. One example are fractals sets, which are often used to model natural phenomena such as plants, clouds, and terrains [\[SKG<sup>+</sup>09\]](#). Fractal objects are generated as the limit of an iterative process [\[Man83\]](#). Due to their regular, recursive structure it is often possible to answer proximity queries about such a set without generating the set itself.

Two other examples of infinite sets generated implicitly from finite models include (1) subdivision surfaces [\[AS10\]](#), where a smooth surface is generated by applying a recursive refinement process to a finite set of boundary points, and (2) metaballs [\[Bli82\]](#), where a surface is defined by a blending function applied to a collection of geometric balls. In both cases, it is possible to answer nearest neighbor queries for the underlying object without the need to generate its boundary.

Proximity queries have been applied before. Panahi *et al.* [\[PASG13\]](#) use proximity probes on a convex polygon in the plane to reconstruct it exactly. Goel *et al.* [\[GIV01\]](#), reduce the approximation versions of several problems like diameter, farthest neighbors, discrete center, metric facility location, bottleneck

matching and minimum weight matching to nearest neighbor queries. They sometimes require other primitives for their algorithms, for example computation of the minimum enclosing ball or a dynamic version of the approximate nearest-neighbor oracle. Similarly, the computation of the minimum spanning tree [HIM12] can be done using nearest-neighbor queries (but the data structure needs to support deletions). For more details, see the survey by Indyk [Ind04].

## Our contributions

In this paper we consider a number of problems on implicitly-represented point sets. Here is a summary of our main results.

**$k$ -center clustering and the greedy permutation.** Given a point set  $P$ , a *greedy permutation* (informally) is an ordering of the points of  $P$ :  $p_1, \dots, p_k, \dots$ , such that for any  $k$ , the set of points  $\{p_1, \dots, p_k\}$  is a  $O(1)$ -approximation to the optimal  $k$ -center clustering. This sequence arises in the  $k$ -center approximation of Gonzalez [Gon85], and its properties were analyzed by Har-Peled and Mendel [HM06]. Specifically, if  $P$  can be covered by  $k$  balls of radius  $r_k$ , then the maximum distance of any point of  $P$  to its nearest neighbor in  $\{p_1, \dots, p_k\}$  is  $O(r_k)$ .

In Section 3, we show that under reasonable assumptions, in constant dimension, one can compute such an approximate greedy permutation using  $O(k)$  exact proximity queries. If the oracle answers  $(1 + \varepsilon)$ -ANN queries, then for any  $k$ , the permutation generated is competitive with the optimal  $k$ -center clustering, considering the first  $O(k \log_{1/\varepsilon} \Phi)$  points in this permutation, where  $\Phi$  is (roughly) the spread of the point set. The hidden constant factors grow exponentially in the dimension.

**Approximate convex-hull membership.** Given a point set  $P$  in  $\mathbb{R}^d$ , we consider the problem of deciding whether a given query point  $q \in \mathbb{R}^d$  is inside its convex-hull  $\mathcal{C} = \mathcal{CH}(P)$ . We say that the answer is  $\varepsilon$ -approximately correct if the answer is correct whenever the query point's distance from the boundary of  $\mathcal{C}$  is at least  $\varepsilon \cdot \text{diam}(\mathcal{C})$ . In Section 4, we show that, given an oracle for  $(1 + \varepsilon^2/c)$ -ANN queries, for some sufficiently large constant  $c$ , it is possible to answer approximate convex-hull membership queries using  $O(1/\varepsilon^2)$  proximity queries. Remarkably, the number of queries is independent of the dimension of the data.

Our algorithm operates iteratively, by employing a gradient descent-like approach. It generates a sequence of points, all within the convex hull, that converges to the query point. Similar techniques have been used before, and are sometimes referred to as the Frank-Wolfe algorithm. Clarkson provides a survey and some new results of this type [Cla10]. A recent algorithm of this type is the work by Kalantari [Kal12]. Our main new contribution for the convex-hull membership problem is showing that the iterative algorithm can be applied to implicit point sets using nearest-neighbor queries.

**Balanced proximity clustering.** We study a problem that involves summarizing a point set in a way that preserves proximity information. Specifically, given a set  $P$  of  $n$  points in  $\mathbb{R}^d$ , and a parameter  $k$ , the objective is to select  $m$  centers from  $P$ , such that if we assign every point of  $P$  to its nearest center, no center has been selected by more than  $k$  points. This problem is related to topic of capacitated clustering from operations research [MB84].

In Section 5, we show that in the plane there exists such a clustering consisting of  $O(n/k)$  such centers, and that in higher dimensions one can select  $O((n/k) \log(n/k))$  centers (where the constant depends on the dimension). This result is not directly related to the other results in the paper.

**Paper organization.** In [Section 2](#) we review some relevant work on  $k$ -center clustering. In [Section 3](#) we provide our algorithm to compute an approximate  $k$ -center clustering. In [Section 4](#) we show how we can decide approximately if a query point is within the convex hull of the given data points in a constant number of queries, where the constant depends on the degree of accuracy desired. Finally, in [Section 5](#) we investigate balanced Voronoi partitions, which provides a density-based clustering of the data. Here we assume that all the data is known and the goal is to come up with a useful clustering that can help in proximity search queries.

## 2. Preliminaries

### 2.1. Background — $k$ -center clustering and the greedy permutation

The following is taken from [[Har11](#), Chap. 4], and is provided here for the sake of completeness.

In the  $k$ -center clustering problem, a set  $P \subseteq \mathbb{R}^d$  of  $n$  points is provided together with a parameter  $k$ . The objective is to find a set of  $k$  points,  $C \subseteq P$ , such that the maximum distance of a point in  $P$  to its closest point in  $C$  is minimized. Formally, define  $\text{price}(C, P) = \max_{p \in P} \min_{c \in C} \|p - c\|$ . Let  $C_{\text{opt}}$  denote the set of centers achieving this minimum. The  $k$ -center problem can be interpreted as the problem of computing the minimum radius, called the  *$k$ -center clustering radius*, such that it is possible to cover the points of  $P$  using  $k$  balls of this radius, each centered at one of the data points. It is known that  $k$ -center clustering is **NP-HARD**. Even in the plane, it is **NP-HARD** to approximate to within a factor of  $(1 + \sqrt{7})/2 \approx 1.82$  [[FG88](#)].

**The greedy clustering algorithm.** Gonzalez [[Gon85](#)] provided a 2-approximation algorithm for  $k$ -center clustering. This algorithm, denoted by **GreedyKCenter**, repeatedly picks the point farthest away from the current set of centers and adds it to this set. Specifically, it starts by picking an arbitrary point,  $\bar{c}_1$ , and setting  $C_1 = \{\bar{c}_1\}$ . For  $i > 1$ , in the  $i$ th iteration, the algorithm computes

$$r_{i-1} = \text{price}(C_{i-1}, P) = \max_{p \in P} d(p, C_{i-1}) \quad (2.1)$$

and the point  $\bar{c}_i$  that realizes it, where  $d(p, C_{i-1}) = \min_{c \in C_{i-1}} \|p - c\|$ . Next, the algorithm adds  $\bar{c}_i$  to  $C_{i-1}$  to form the new set  $C_i$ . This process is repeated until  $k$  points have been collected.

If we run **GreedyKCenter** till it exhausts all the points of  $P$  (i.e.,  $k = n$ ), then this algorithm generates a permutation of  $P$ ; that is,  $\langle P \rangle = \langle \bar{c}_1, \dots, \bar{c}_n \rangle$ . We will refer to  $\langle P \rangle$  as the *greedy permutation* of  $P$ . There is also an associated sequence of radii  $\langle r_1, \dots, r_n \rangle$ , and the key property of the greedy permutation is that for each  $i$  with  $1 \leq i \leq n$ , all the points of  $P$  are within a distance at most  $r_i$  from the points of  $C_i = \langle \bar{c}_1, \dots, \bar{c}_i \rangle$ . The greedy permutation has applications to packings, which we describe next.

**Definition 2.1.** A set  $S \subseteq P$  is an  *$r$ -packing* for  $P$  if the following two properties hold:

- (i) **Covering property**: All the points of  $P$  are within a distance at most  $r$  from the points of  $S$ .
- (ii) **Separation property**: For any pair of points  $p, x \in S$ ,  $\|p - x\| \geq r$ .

(For most purposes, one can relax the separation property by requiring that the points of  $S$  be at distance  $\Omega(r)$  from each other.)

Intuitively, an  $r$ -packing of a point set  $P$  is a compact representation of  $P$  at resolution  $r$ . Surprisingly, the greedy permutation of  $P$  provides us with such a representation for all resolutions.

**Lemma 2.2** ([Har11]). (A) Let  $P$  be a set of  $n$  points in  $\mathbb{R}^d$ , and let its greedy permutation be  $\langle \bar{c}_1, \dots, \bar{c}_n \rangle$  with the associated sequence of radii  $\langle r_1, \dots, r_n \rangle$ . For any  $i$ ,  $C_i = \langle \bar{c}_1, \dots, \bar{c}_i \rangle$  is an  $r_i$ -packing of  $P$ . Furthermore,  $r_i$  is a 2-approximation for the optimal  $i$ -center clustering radius of  $P$ .

(B) For any  $k$ , let  $r_{\text{opt}}^k$  be the radius of the optimal  $k$ -center clustering of  $P$ . Then, for any constant  $c$ ,  $r_{\text{opt}}^{O(c^d k)} \leq r_{\text{opt}}^k / c$ .

(C) Computing the optimal  $k$ -center clustering of the first  $O(k/\varepsilon^d)$  points of the greedy permutation, after appropriate rescaling, results in a  $(1 + \varepsilon)$ -approximation to the optimal  $k$ -center clustering of  $P$ .

## 2.2. Setup

Our algorithms operate on a (not necessarily finite) point set  $P$  in  $\mathbb{R}^d$ . We assume that we are given a compact subset of  $\mathbb{R}^d$ , called the **domain** and denoted  $\mathcal{D}$ , such that  $P \subseteq \mathcal{D}$ . Throughout we assume that  $\mathcal{D}$  is the unit hypercube  $[0, 1]^d$ . The set  $P$  (not necessarily finite) is contained in  $\mathcal{D}$ .

Given a query point  $q \in [0, 1]^d$ , let  $\text{nn}(q, P) = \arg \min_{p \in P} \|q - p\|$  denote the nearest neighbor (NN) of  $q$ . We say a point  $x$  is a  $(1 + \varepsilon)$ -approximate nearest-neighbor (ANN) for  $q$  if  $\|q - x\| \leq (1 + \varepsilon) \|q - \text{nn}(q, P)\|$ . We assume that the sole access to  $P$  is through “black-box” data structures  $T_{nn}$  and  $T_{ann}$ , which given a query point  $q$ , return the NN and ANN, respectively, to  $q$  in  $P$ .

## 3. Using proximity search to compute $k$ -center clustering

**The problem.** Our purpose is to compute (or approximately compute) a  $k$ -center clustering of  $P$  through the ANN black box we have, where  $k$  is a given parameter between 1 and  $n$ .

### 3.1. Greedy permutation via NN queries: **GreedyPermutNN**

Let  $q_0$  be an arbitrary point in  $\mathcal{D}$ . Let  $\nu_0$  be its nearest-neighbor in  $P$  computed using the provided NN data structure  $T_{nn}$ . Let  $b_0 = \text{ball}(q_0, \|q_0 - \nu_0\|)$  be the open ball of radius  $\|q_0 - \nu_0\|$  centered at  $q_0$ . Finally, let  $G_0 = \{\nu_0\}$ , and let  $\mathcal{D}_0 = \mathcal{D} \setminus b_0$ .

In the  $i$ th iteration, for  $i > 0$ , let  $q_i$  be the point in  $\mathcal{D}_{i-1}$  farthest away from  $G_{i-1}$ . Formally, this is the point in  $\mathcal{D}_{i-1}$  that maximizes  $d(q_i, G_{i-1})$ , where  $d(q, X) = \max_{c \in X} \|q - c\|$ . Let  $\nu_i = \text{nn}(q_i, P)$  denote the nearest-neighbor  $\nu_i$  to  $q_i$  in  $P$ , computed using  $T_{nn}$ . Let

$$r_i = d(q_i, G_{i-1}), \quad b_i = \text{ball}(q_i, r_i), \quad G_i = G_{i-1} \cup \{\nu_i\}, \quad \text{and} \quad \mathcal{D}_i = \mathcal{D}_{i-1} \setminus b_i.$$

Left to its own devices, this algorithm computes a sequence of not necessarily distinct points  $\nu_0, \nu_1, \dots$  of  $P$ . If  $P$  is not finite then this sequence may also have infinitely many distinct points. Furthermore,  $\mathcal{D}_0 \supseteq \mathcal{D}_1 \supseteq \dots$  is a sequence of outer approximations to  $P$ .

The execution of this algorithm is illustrated in **Figure 3.1**.

### 3.2. Analysis

Let  $\mathcal{O} = \{o_1, \dots, o_k\}$  be an optimal set of  $k$  centers of  $P$ . Formally, it is a set of  $k$  points in  $P$  that minimizes the quantity  $r_{\text{opt}}^k = \max_{q \in P} d(q, \mathcal{O})$ . Specifically,  $r_{\text{opt}}^k$  is the smallest possible radius such that  $k$  closed balls of that radius centered at points in  $P$ , cover  $P$ . Our claim is that after  $O(k)$  iterations of the algorithm **GreedyPermutNN**, the sequence of points provides a similar quality clustering of  $P$ .

For any given point  $p \in \mathbb{R}^d$  we can cover the sphere of directions centered at  $p$  by narrow cones of angular diameter at most  $\pi/12$ . We fix such a covering, denoting the set of cones by  $\mathcal{C}_p$ , and observe

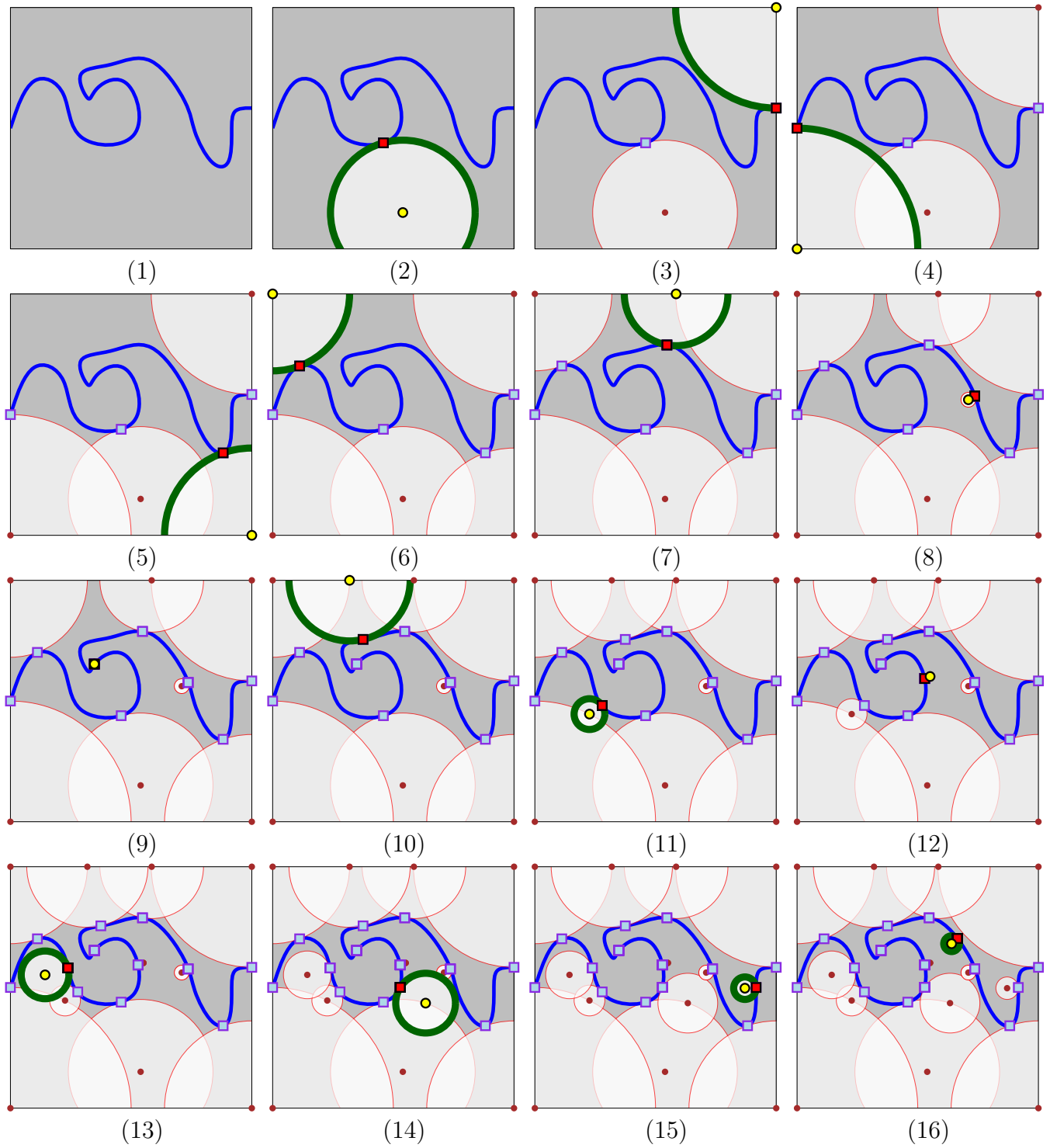


Figure 3.1: An example of the execution of the algorithm **GreedyPermutNN** of Section 3.1.



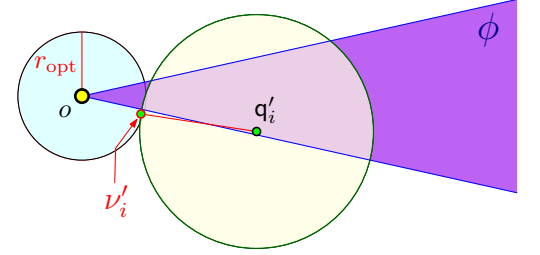
that the number of such cones is a constant  $c_d$  that depends on the dimension. Moreover, by simple translation we can transfer such a covering to be centered at any point  $\mathbf{p}' \in \mathbb{R}^d$ .

**Lemma 3.1.** *After  $\mu = kc_d$  iterations, for any optimal center  $o_i \in \mathcal{O}$ , we have  $d(o_i, G_\mu) \leq 3r_{\text{opt}}$ , where  $r_{\text{opt}} = r_{\text{opt}}^k$ .*

*Proof:* If for any  $j \leq \mu$ , we have  $r_j \leq 3r_{\text{opt}}$  then all the points of  $\mathcal{D}_{j-1} \supseteq \mathbf{P}$  are in distance at most  $3r_{\text{opt}}$  from  $G_j$ , and the claim trivially holds as  $\mathcal{O} \subseteq \mathbf{P}$ .

Let  $o$  be an optimal center and let  $\mathbf{P}_o$  be the set of points of  $\mathbf{P}$  that are closest to  $o$  among all the centers of  $\mathcal{O}$ , i.e.,  $\mathbf{P}_o$  is the cluster of  $o$  in the optimal clustering. Fix a cone  $\phi$  from  $\mathcal{C}_o$  ( $\phi$ 's apex is at  $o$ ). Consider the output sequence  $\nu_0, \nu_1, \dots$ , and the corresponding query sequence  $\mathbf{q}_0, \mathbf{q}_1, \dots$  computed by the algorithm. In the following, we use the property of the algorithm that  $r_1 \geq r_2 \geq \dots$ , where  $r_i = d(\mathbf{q}_i, G_{i-1})$ . A point  $\mathbf{q}_j$  is **admissible** if (i)  $\nu_j \in \mathbf{P}_o$ , and (ii)  $\mathbf{q}_j \in \phi$  (in particular,  $\nu_j$  is not necessarily in  $\phi$ ).

We proceed to show that there are at most  $O(1)$  admissible points for a fixed cone, which by a packing argument will imply the claim as every  $\mathbf{q}_j$  is admissible for exactly one cone. Consider the induced subsequence of the output sequence restricted to the admissible points of  $\phi$ :  $\nu'_1, \nu'_2, \dots$ , and let  $\mathbf{q}'_1, \mathbf{q}'_2, \dots$  be the corresponding query points used by the algorithm. Formally, for a point  $\nu'_i$  in this sequence, let  $\text{iter}(i)$  be the iteration of the algorithm it was created. Thus, for all  $i$ , we have  $\mathbf{q}'_i = \mathbf{q}_{\text{iter}(i)}$  and  $\nu'_i = \nu_{\text{iter}(i)}$ .



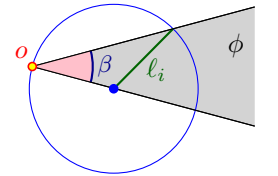
Observe that  $\mathbf{P}_o \subseteq \mathbf{P} \cap \text{ball}(o, r_{\text{opt}})$ . This implies that

$$\|\nu'_j - o\| \leq r_{\text{opt}}, \quad \text{for all } j.$$

Let  $\ell'_i = \|\mathbf{q}'_i - \nu'_i\|$  and  $r'_i = d(\mathbf{q}'_i, G_{\text{iter}(i)-1})$ . Observe that for  $i > 1$ , we have  $\ell'_i \leq r'_i \leq \ell'_i + 2r_{\text{opt}}$ , as  $\nu'_{i-1} \in \mathbf{P}_o$ . Hence, if  $\ell'_i \leq r_{\text{opt}}$ , then  $r'_i \leq 3r_{\text{opt}}$ , and we are done. This implies that for any  $i, j$ , such that  $1 < i < j$ , it must be that  $\|\mathbf{q}'_i - \mathbf{q}'_j\| \geq \ell'_i > r_{\text{opt}}$ , as the algorithm carves out a ball of radius  $\ell'_i$  around  $\mathbf{q}'_i$ , and  $\mathbf{q}'_j$  must be outside this ball.

By a standard packing argument, there can be only  $O(1)$  points in the sequence  $\mathbf{q}'_2, \mathbf{q}'_3, \dots$  that are within distance at most  $10r_{\text{opt}}$  from  $o$ . If there are no points beyond this distance, we are done. Otherwise, let  $i > 1$  be the minimum index, such that  $\mathbf{q}'_i$  is at distance larger than  $10r_{\text{opt}}$  from  $o$ . We now prove that the points of  $\phi \setminus \text{ball}(\mathbf{q}'_i, \ell'_i)$  are of two types — those contained within  $\text{ball}(o, 3r_{\text{opt}})$  and those that lie at distance greater than  $(4/3)\ell'_i$  from  $o$ .

To see this, observe that since the angle of the cone was chosen to be sufficiently small,  $\text{ball}(\mathbf{q}'_i, \ell'_i)$  splits  $\phi$  into two components, where all the points in the component containing  $o$  are distance  $< 3r_{\text{opt}}$  from  $o$ . The minimum distance to  $o$  (from a point in the component not containing  $o$ ) is realized when  $\mathbf{q}'_i$  is on the boundary of  $\phi$  and  $o$  is on the boundary of  $\text{ball}(\mathbf{q}'_i, \ell'_i)$ . Then the distance of any point of  $\phi \setminus \text{ball}(\mathbf{q}'_i, \ell'_i)$  from  $o$  is at least  $2\ell'_i \cos(\beta) \geq 2\ell'_i \sqrt{3/4} \geq 1.73\ell'_i$ , as the opening angle of the cone is at most  $\pi/12$ . See figure on the right. The general case is somewhat more complicated as  $o$  might be in distance at most  $r_{\text{opt}}$  from the boundary of  $\text{ball}(\mathbf{q}'_i, \ell'_i)$ , but as  $\ell'_i \geq 10r_{\text{opt}}$ , the claim still holds — we omit the tedious but straightforward calculations.



In particular, this implies that any later point  $\mathbf{q}'_k$  in the sequence (i.e.,  $k > i$ ) is either one of the  $O(1)$  close points, or it must be far away, but then it is easy to argue that  $r'_k$  must be larger than  $r'_i$ , which is a contradiction as  $r_2 \geq r_3 \geq \dots$  (as  $r'_i$  appears before  $r'_k$  in this sequence). ■

The above lemma readily implies the following.

**Theorem 3.2.** *Let  $P \subseteq \mathcal{D}$  be a given set of points in  $\mathbb{R}^d$  (not necessarily finite), where  $\mathcal{D}$  is a bounded set in  $\mathbb{R}^d$ . Furthermore, assume that  $P$  can be accessed only via a data structure  $T_{nn}$  that answers exact nearest-neighbor (NN) queries on  $P$ . The algorithm **GreedyPermutNN**, described in [Section 3.1](#), computes a permutation  $\langle \nu_0, \dots \rangle$  of  $P$ , such that, for any  $k > 0$ ,  $P \subseteq \bigcup_{i=0}^{ck} \text{ball}(\nu_i, r_{\text{opt}}^k)$ , where  $c$  is a constant (independent of  $k$ ), and  $r_{\text{opt}}^k$  is the minimum radius of  $k$  balls (of the same radius) needed to cover  $P$ .*

*The algorithm can be implemented, such that running it for  $i$  iterations, takes polynomial time in  $i$  and involves  $i$  calls to  $T_{nn}$ .*

*Proof:* Using [Lemma 2.2 \(B\)](#) in [Lemma 3.1](#) implies the result. As for the running time, naively one needs to maintain the arrangement of balls inside the domain, and this can be done in polynomial time in the number of balls. ■

**Observation 3.3.** *If  $P$  is finite of size  $n$ , the above theorem implies that after  $i \geq cn$  iterations, one can recover the entire point set  $P$  (as  $r_{\text{opt}}^n = 0$ ). Therefore  $cn$  is an upper bound on the number of queries for any problem. Note however that in general our goal is to demonstrate when problems can be solved using a significantly smaller amount of NN queries.*

The above also implies an algorithm for approximating the diameter.

**Lemma 3.4.** *Consider the setting of [Theorem 3.2](#) using an exact nearest-neighbor oracle. Suppose that the algorithm is run for  $m = c_d + 1$  iterations, and let  $\nu_1, \dots, \nu_m$  be the set of output centers and  $r_1, \dots, r_m$  be the corresponding distances. Then,  $\text{diam}(P)/3 \leq \max(\text{diam}(\nu_1, \dots, \nu_m), r_m) \leq 3 \cdot \text{diam}(P)$ .*

*Proof:* Since the discrete one-center clustering radius lies in the interval  $[\text{diam}(P)/2, \text{diam}(P)]$ , [Lemma 3.1](#) implies that  $r_m \leq 3r_{\text{opt}} \leq 3 \cdot \text{diam}(P)$ . Moreover, each  $\nu_i$  is in  $P$ , and so  $\text{diam}(\nu_1, \dots, \nu_m) \leq \text{diam}(P)$ . Thus the upper bound follows.

For the lower bound, observe that if  $\text{diam}(\nu_1, \dots, \nu_m) < \text{diam}(P)/3$ , as well as  $r_m < \text{diam}(P)/3$ , then it must be true that  $P \subseteq \mathcal{D}_{m-1} \subseteq \bigcup_{j=1}^l \text{ball}(\nu_j, r_m)$  has diameter less than  $\text{diam}(P)$ , a contradiction. ■

### 3.3. Using approximate nearest-neighbor search

If we are using an ANN black box  $T_{ann}$  to implement the algorithm, one can no longer scoop away the ball  $b_i = \text{ball}(\mathbf{q}_i, \|\mathbf{q}_i - \nu_i\|)$  at the  $i$ th iteration, as it might contain some of the points of  $P$ . Instead, one has to be more conservative, and use the ball  $b'_i = \text{ball}(\mathbf{q}_i, (1 - \varepsilon) \|\mathbf{q}_i - \nu_i\|)$ . Now, we might need to perform several queries till the volume being scooped away is equivalent to a single exact query.

Specifically, let  $P$  be a finite set, and consider its associated **spread**:

$$\Phi = \frac{\text{diam}(\mathcal{D}_0)}{\min_{\mathbf{p}, \mathbf{x} \in P} \|\mathbf{p} - \mathbf{x}\|}.$$

We can no longer claim, as in [Lemma 3.1](#), that each cone would be visited only one time (or constant number of times). Instead, it is easy to verify that each query point in the cone, shrinks the diameter of the domain restricted to the cone by a factor of roughly  $\varepsilon$ . As such, at most  $O(\log_{1/\varepsilon} \Phi) = O(\varepsilon^{-1} \log \Phi)$  query points would be associated with each cone.

**Corollary 3.5.** *Consider the setting of [Theorem 3.2](#), with the modification that we use a  $(1 + \varepsilon)$ -ANN data structure  $T_{ann}$  to access  $P$ . Then, for any  $k$ ,  $P \subseteq \bigcup_{i=1}^{f(k)} \text{ball}(\nu_i, r_{\text{opt}}^k)$ , where  $f(k) = O(k\varepsilon^{-1} \log \Phi)$ .*



### 3.4. Discussion

**Outer approximation.** As implied by the algorithm description, one can think about the algorithm providing an outer approximation to the set:  $\mathcal{D}_1 \supseteq \mathcal{D}_2 \supseteq \dots \supseteq \mathcal{P}$ . As demonstrated in Figure 3.1, the sequence of points computed by the algorithm seems to be a reasonable greedy permutation of the underlying set. However, the generated outer approximation seems to be inferior. If the purpose is to obtain a better outer approximation, a better strategy may be to pick the  $i$ th query point  $\mathbf{q}_i$  as the point inside  $\mathcal{D}_i$  farthest away from  $\partial\mathcal{D}_{i-1} \cup G_{i-1}$ .

**Implementation details.** We have not spent any effort to describe in detail the algorithm of Theorem 3.2, mainly because an implementation of the exact version seems quite challenging in practice. A more practical approach would be to describe the uncovered domain  $\mathcal{D}_i$  approximately, by approximating from the inside, every ball  $b_i$  by an  $O(1/\varepsilon^d)$  grid of cubes, and maintaining these cubes using a (compressed) quadtree. This provides an explicit representation of the complement of the union of the approximate balls. Next, one would need to maintain for every free leaf of this quadtree, a list of points of  $G_i$  that might serve as its nearest neighbors — in the spirit of approximate Voronoi diagrams [Har11].

## 4. Convex-hull membership queries via proximity queries

Let  $\mathcal{P}$  be a set of  $n$  points in  $\mathbb{R}^d$ , let  $\Delta$  denote  $\mathcal{P}$ 's diameter, and let  $\varepsilon > 0$  be a prespecified parameter. We assume that the value of  $\Delta$  is known, although a constant approximation to this value is sufficient for our purposes. (See Lemma 3.4 on how to compute this under reasonable assumptions.)

Let  $\mathcal{C} = \mathcal{CH}(\mathcal{P})$  denote  $\mathcal{P}$ 's convex hull. Given a query point  $\mathbf{q} \in \mathbb{R}^d$ , the task at hand is to decide if  $\mathbf{q}$  is in  $\mathcal{C}$ . As before, we assume that our only access to  $\mathcal{P}$  is via an ANN data structure. There are two possible outputs:

- (A) IN: if  $\mathbf{q} \in \mathcal{C}$ , and
- (B) OUT: if  $\mathbf{q}$  is at distance greater than  $\varepsilon\Delta$  from  $\mathcal{C}$ ,

Either answer is acceptable if  $\mathbf{q}$  lies within distance  $\varepsilon\Delta$  of  $\partial\mathcal{C}$ .

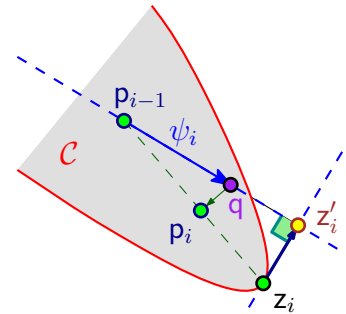
### 4.1. Convex hull membership queries using exact extremal queries

We first solve the problem using exact extremal queries and then later show these queries can be answered approximately with ANN queries.

#### 4.1.1. The algorithm

We construct a sequence of points  $\mathbf{p}_0, \mathbf{p}_1, \dots$  each guaranteed to be in the convex hull  $\mathcal{C}$  of  $\mathcal{P}$  and use them to determine whether  $\mathbf{q} \in \mathcal{C}$ . The algorithm is as follows. Let  $\mathbf{p}_0$  be an arbitrary point of  $\mathcal{P}$ . For  $i > 0$ , in the  $i$ th iteration, the algorithm checks whether  $\|\mathbf{p}_{i-1} - \mathbf{q}\| \leq \varepsilon\Delta$ , and if so the algorithm outputs IN and stops.

Otherwise, consider the ray  $\psi_i$  emanating from  $\mathbf{p}_{i-1}$  in the direction of  $\mathbf{q}$ . The algorithm computes the point  $\mathbf{z}_i \in \mathcal{P}$  that is extremal in the direction of this ray. If the projection  $\mathbf{z}'_i$  of  $\mathbf{z}_i$  on the line supporting  $\psi_i$  is between  $\mathbf{p}_{i-1}$  and  $\mathbf{q}$ , then  $\mathbf{q}$  is outside the convex-hull  $\mathcal{C}$ , and the algorithm stops and returns OUT. Otherwise, the algorithm sets  $\mathbf{p}_i$



to be the projection of  $\mathbf{q}$  on the line segment  $\mathbf{p}_{i-1}\mathbf{z}_i$ , and continues to the next iteration. See figure on the right, and [Figure 4.1](#).

The algorithm performs  $O(1/\varepsilon^2)$  iterations, and returns OUT if it did not stop earlier.

#### 4.1.2. Analysis

**Lemma 4.1.** *If the algorithm runs for more than  $i$  iterations, then  $\mathbf{d}_i < \left(1 - \frac{\varepsilon^2}{2}\right)\mathbf{d}_{i-1}$ , where  $\mathbf{d}_i = \|\mathbf{q} - \mathbf{p}_i\|$ .*

*Proof:* By construction,  $\mathbf{p}_i$ ,  $\mathbf{p}_{i-1}$ , and  $\mathbf{q}$  form a right angle triangle. The proof now follows by a direct trigonometric argument. Consider [Figure 4.1](#). We have the following properties:

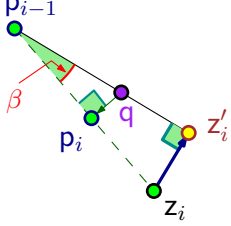


Figure 4.1

- (A) The triangles  $\triangle \mathbf{p}_{i-1}\mathbf{z}'_i\mathbf{z}_i$  and  $\triangle \mathbf{p}_{i-1}\mathbf{p}_i\mathbf{q}$  are similar.
- (B) Because the algorithm has not terminated in the  $i$ th iteration,  $\|\mathbf{p}_{i-1} - \mathbf{q}\| > \varepsilon\Delta$ .
- (C) The point  $\mathbf{q}$  must be between  $\mathbf{p}_{i-1}$  and  $\mathbf{z}'_i$ , as otherwise the algorithm would have terminated. Thus,  $\|\mathbf{p}_{i-1} - \mathbf{z}'_i\| \geq \|\mathbf{p}_{i-1} - \mathbf{q}\| > \varepsilon\Delta$ .
- (D) We have  $\|\mathbf{p}_{i-1} - \mathbf{z}_i\| \leq \Delta$ , since both points are in  $\mathcal{C}$ .

We conclude that  $\cos \beta = \frac{\|\mathbf{p}_{i-1} - \mathbf{z}'_i\|}{\|\mathbf{p}_{i-1} - \mathbf{z}_i\|} > \frac{\varepsilon\Delta}{\Delta} = \varepsilon$ . Now, we have

$$\begin{aligned} \|\mathbf{q} - \mathbf{p}_i\| &= \|\mathbf{q} - \mathbf{p}_{i-1}\| \sin \beta = \|\mathbf{q} - \mathbf{p}_{i-1}\| \sqrt{1 - \cos^2 \beta} < \sqrt{1 - \varepsilon^2} \|\mathbf{q} - \mathbf{p}_{i-1}\| \\ &< \left(1 - \frac{\varepsilon^2}{2}\right) \|\mathbf{q} - \mathbf{p}_{i-1}\|, \end{aligned}$$

since  $(1 - \varepsilon^2/2)^2 > 1 - \varepsilon^2$ . ■

**Lemma 4.2.** *Either the algorithm stops within  $O(1/\varepsilon^2)$  iterations with a correct answer, or the query point is more than  $\varepsilon\Delta$  far from the convex hull  $\mathcal{C}$ ; in the latter case, since the algorithm says OUT its output is correct.*

*Proof:* If the algorithm stops before it completes the maximum number of iterations, it can be verified that the output is correct as there is an easy certificate for this in each of the possible cases.

Otherwise, suppose that the query point is within  $\varepsilon\Delta$  of  $\mathcal{C}$ . We argue that this leads to a contradiction; thus the query point must be more than  $\varepsilon\Delta$  far from  $\mathcal{C}$  and the output of the algorithm is correct. Observe that  $\mathbf{d}_i$  is a monotone decreasing quantity that starts at values  $\leq \Delta$  (i.e,  $\mathbf{d}_0 \leq \Delta$ ), since otherwise the algorithm terminates after the first iteration, as  $\mathbf{z}'_1$  would be between  $\mathbf{q}$  and  $\mathbf{p}_0$  on  $\psi_1$ .

Consider the  $j$ th **epoch** to be block of iterations of the algorithm, where  $2^{-j}\Delta < \mathbf{d}_i \leq 2^{-j+1}\Delta$ . Following the proof of [Lemma 4.1](#), one observes that during the  $j$ th epoch one can set  $\varepsilon_j = 1/2^j$  in place of  $\varepsilon$ , and using the argument it is easy to show that the  $j$ th epoch lasts  $O(1/\varepsilon_j^2)$  iterations. By assumption, since the algorithm continued for the maximum number of iterations we have  $\mathbf{d}_i > \varepsilon\Delta$ , and so the maximum number of epochs is  $\lceil \lg(1/\varepsilon) \rceil$ . As such, the total number of iterations is  $\sum_{j=1}^{\lceil \lg(1/\varepsilon) \rceil} O(1/\varepsilon_j^2) = O(1/\varepsilon^2)$ . Since the algorithm did not stop this is a contradiction. ■

### 4.1.3. Approximate extremal queries

For our purposes, approximate extremal queries on  $\mathbf{P}$  are sufficient.

**Definition 4.3.** A data structure provides  *$\epsilon$ -approximate extremal queries* for  $P$ , if for any query unit vector  $\mathbf{v}$ , it returns a point  $\mathbf{p}$ , such that

$$\forall x \in P, \quad \langle v, x \rangle \leq \langle v, p \rangle + \varepsilon \cdot \text{diam}(P),$$

where  $\langle \mathbf{v}, \mathbf{x} \rangle$  denotes the dot-product of  $\mathbf{v}$  with  $\mathbf{x}$ .

One can now modify the algorithm of [Section 4.1.1](#) to use, say,  $\varepsilon/4$ -approximate extremal queries on  $P$ . Indeed, one modifies the algorithm so it stops only if  $\mathbf{z}_i$  is on the segment  $\mathbf{p}_{i-1}\mathbf{q}$ , and it is in distance more than  $\varepsilon\Delta/4$  away from  $\mathbf{q}$ . Otherwise the algorithm continues. It is straightforward but tedious to prove that the same algorithm performs asymptotically the same number of iterations (intuitively, all that happens is that the constants get slightly worse). The worse case as far progress in a single iteration is depicted in [Figure 4.2](#).

**Lemma 4.4.** *The algorithm of Section 4.1.1 can be modified to use  $\varepsilon/4$ -approximate extremal queries and output a correct answer after performing  $O(1/\varepsilon^2)$  iterations.*

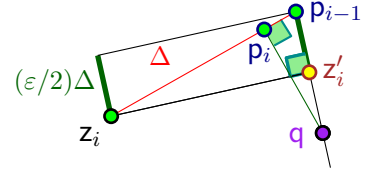


Figure 4.2: Worse case if extremal queries are approximate.

## 4.2. Convex-hull membership via ANN queries

#### 4.2.1. Approximate extremal queries via ANN queries

The basic idea is to replace the extremal empty half-space query, by an ANN query. Specifically, a  $(1 + \delta)$ -ANN query performed at  $\mathbf{q}$  returns us a point  $\mathbf{p}$ , such that

$$\forall \mathbf{x} \in \mathbf{P}, \quad \|\mathbf{q} - \mathbf{p}\| \leq (1 + \delta) \|\mathbf{q} - \mathbf{x}\|.$$

Namely,  $\text{ball}\left(\mathbf{q}, \frac{\|\mathbf{q}-\mathbf{p}\|}{1+\delta}\right)$  does not contain any points of  $\mathbf{P}$ . Locally, a ball looks like a halfspace, and so by taking the query point to be sufficiently far and the approximation parameter to be sufficiently small, the resulting empty ball and its associated **ANN** can be used as the answer to an extremal direction query.

### 4.2.2. The modified algorithm

Assume the algorithm is given a data structure  $T_{ann}$  that can answer  $(1 + \delta)$ -ANN queries on  $\mathbf{P}$ . Also assume that it is provided with an initial point  $\mathbf{p}_0 \in \mathbf{P}$ , and a value  $\Delta'$  that is, say, a 2-approximation to  $\Delta = \text{diam}(\mathbf{P})$ , that is  $\Delta \leq \Delta' \leq 2\Delta$ .

In the  $i$ th iteration, the algorithm considers (again) the ray  $\psi_i$  starting from  $\mathbf{p}_i$ , in the direction of  $\mathbf{q}$ . Let  $\mathbf{q}_i$  be the point within distance, say,

$$\tau = c\Delta'/\varepsilon \quad (4.1)$$

from  $\mathbf{p}_{i-1}$  along  $\psi_i$ , where  $c$  is an appropriate constant to be determined shortly. Next, let  $\mathbf{z}_i$  be the  $(1 + \delta)$ -ANN returned by  $T_{ann}$  for the query point  $\mathbf{q}_i$ , where the value of  $\delta$  would be specified shortly. The algorithm now continues as before, by setting  $\mathbf{p}_i$  to be the nearest point on  $\mathbf{p}_{i-1}\mathbf{z}_i$  to  $\mathbf{q}$ . Naturally, if  $\|\mathbf{q} - \mathbf{p}_i\|$  falls below  $\varepsilon\Delta'/2$ , the algorithm stops, and returns IN, and otherwise the algorithm continues to the next iteration. As before, if the algorithm reaches the  $N$ th iteration, it stops and returns OUT, where  $N = O(1/\varepsilon^2)$ .

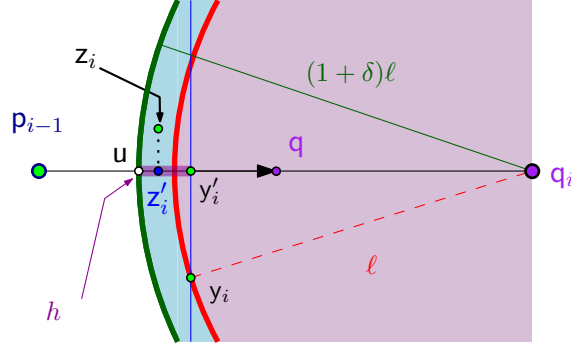


Figure 4.3: Illustration of the proof of Lemma 4.5.

### 4.2.3. Analysis

**Lemma 4.5.** *Let  $0 < \varepsilon \leq 1$  be a prespecified parameter, and let  $\delta = \varepsilon^2/(32 - \varepsilon)^2 = O(\varepsilon^2)$ . Then, a  $(1 + \delta)$ -ANN query done using  $\mathbf{q}_i$  performed by the algorithm, returns a point  $\mathbf{z}_i$  which is a valid  $\varepsilon$ -approximate extremal query on  $\mathbf{P}$ , in the direction of  $\psi_i$ .*

*Proof:* Consider the extreme point  $\mathbf{y}_i \in \mathbf{P}$  in the direction of  $\psi_i$ . Let  $\mathbf{y}'_i$  be the projection of  $\mathbf{y}_i$  to the segment  $\mathbf{p}_{i-1}\mathbf{q}_i$ , and let  $\ell = \|\mathbf{q}_i - \mathbf{y}_i\|$ . See Figure 4.3.

The  $(1 + \delta)$ -ANN to  $\mathbf{q}_i$  (i.e., the point  $\mathbf{z}_i$ ), must be inside the ball  $b = \text{ball}(\mathbf{q}_i, (1 + \delta)\ell)$ , and let  $\mathbf{z}'_i$  be its projection to the segment  $\mathbf{p}_{i-1}\mathbf{q}_i$ .

Now, if we interpret  $\mathbf{z}_i$  as the returned answer for the approximate extremal query, then the error is the distance  $\|\mathbf{z}'_i - \mathbf{y}'_i\|$ , which is maximized if  $\mathbf{z}'_i$  is as close to  $\mathbf{p}_{i-1}$  as possible. In particular, let  $\mathbf{u}$  be the point in distance  $(1 + \delta)\ell$  from  $\mathbf{q}_i$  along the segment  $\mathbf{p}_{i-1}\mathbf{q}_i$ . We then have that  $\|\mathbf{z}'_i - \mathbf{y}'_i\| \leq h = \|\mathbf{u} - \mathbf{y}'_i\|$ . Now, since  $\|\mathbf{y}'_i - \mathbf{y}_i\| \leq \|\mathbf{p}_{i-1} - \mathbf{y}_i\| \leq \Delta'$ , we have

$$\begin{aligned} h = \|\mathbf{u} - \mathbf{y}'_i\| &\leq (1 + \delta)\ell - \|\mathbf{y}'_i - \mathbf{q}_i\| = (1 + \delta)\ell - \sqrt{\ell^2 - \|\mathbf{y}'_i - \mathbf{y}_i\|^2} \\ &\leq (1 + \delta)\ell - \sqrt{\ell^2 - (\Delta')^2} = \frac{(1 + \delta)^2\ell^2 - \ell^2 + (\Delta')^2}{(1 + \delta)\ell + \sqrt{\ell^2 - (\Delta')^2}} \leq \frac{(2\delta + \delta^2)\ell^2 + (\sqrt{\delta}\ell)^2}{\ell} \\ &\leq \frac{4\delta\ell^2}{\ell} = 4\delta\ell, \end{aligned}$$

since  $\delta \leq 1$ , and assuming that  $\Delta' \leq \sqrt{\delta}\ell$ . For our purposes, we need that  $4\delta\ell \leq \varepsilon\Delta$ . Both of these constraints translate to the inequalities,  $\left(\frac{\Delta'}{\ell}\right)^2 \leq \delta \leq \frac{\varepsilon\Delta}{4\ell}$ . Observe that, by the triangle inequality, it follows that

$$\ell = \|\mathbf{q}_i - \mathbf{y}_i\| \leq \|\mathbf{q}_i - \mathbf{p}_{i-1}\| + \|\mathbf{p}_{i-1} - \mathbf{y}_i\| \leq \tau + \Delta.$$

A similar argument implies that  $\ell \geq \tau - \Delta$ . In particular, it is enough to satisfy the constraint  $\left(\frac{\Delta'}{\tau - \Delta}\right)^2 \leq \delta \leq \frac{\varepsilon\Delta}{4(\tau + \Delta)}$ , which is satisfied if  $\left(\frac{\Delta'}{\tau - \Delta}\right)^2 \leq \delta \leq \frac{\varepsilon\Delta'/2}{4(\tau + \Delta')}$ , as  $\Delta \leq \Delta' \leq 2\Delta$ . Substituting the value of  $\tau = c\Delta'/\varepsilon$ , see Eq. (4.1), this is equivalent to  $\left(\frac{1}{c/\varepsilon - 1}\right)^2 \leq \delta \leq \frac{\varepsilon/2}{4(c/\varepsilon + 1)}$ , which holds for  $c = 32$ , as can be easily verified, and setting  $\delta = \varepsilon^2/(32 - \varepsilon)^2 = O(\varepsilon^2)$ .  $\blacksquare$

**Theorem 4.6.** *Given a set  $P$  of  $n$  points in  $\mathbb{R}^d$ , let  $\varepsilon \in (0, 1]$  be a parameter, and let  $\Delta'$  be a constant approximation to the diameter of  $P$ . Assume that you are given a data structure that can answer  $(1 + \delta)$ -ANN queries on  $P$ , for  $\delta = O(\varepsilon^2)$ . Then, given a query point  $q$ , one can decide, by performing  $O(1/\varepsilon^2)$   $(1 + \delta)$ -ANN queries whether  $q$  is inside the convex-hull  $\mathcal{C} = \mathcal{CH}(P)$ . Specifically, the algorithm returns*

- IN: if  $q \in \mathcal{C}$ , and
- OUT: if  $q$  is more than  $\varepsilon\Delta$  away from  $\mathcal{C}$ , where  $\Delta = \text{diam}(P)$ .

*The algorithm is allowed to return either answer if  $q \notin \mathcal{C}$ , but the distance of  $q$  from  $\mathcal{C}$  is at most  $\varepsilon\Delta$ .*

## 5. Density clustering

### 5.1. Definition

Given a set  $P$  of  $n$  points in  $\mathbb{R}^d$ , and a parameter  $k$  with  $1 \leq k \leq n$ , we are interested in computing a set  $C \subseteq P$  of “centers”, such that each center is assigned at most  $k$  points, and the number of centers is (roughly)  $n/k$ . In addition, we require that:

- (A) A point of  $P$  is assigned to its nearest neighbor in  $C$  (i.e.,  $C$  induces a *Voronoi partition* of  $P$ ).
- (B) The centers come from the original point set.

Intuitively, this clustering tries to capture the local density — in areas where the density is low, the clusters can be quite large (in the volume they occupy), but in regions with high density the clusters have to be tight and relatively “small”.

Formally, given a set of centers  $C$ , and a center  $c \in C$ , its *cluster* is

$$P_c = \left\{ p \in P \mid \|c - p\| < d(p, C \setminus \{c\}) \right\},$$

where  $d(c, X) = \min_{p \in X} \|c - p\|$  (and assuming for the sake of simplicity of exposition that all distances are distinct). The resulting *clustering* is  $\Pi(P, C) = \{P_c \mid c \in C\}$ . A set of points  $P$ , and a set of centers  $C \subseteq P$  is a *k-density clustering* of  $P$  if for any  $c \in C$ , we have  $|P_c| \leq k$ . As mentioned, we want to compute a balanced partitioning, i.e., one where the number of centers is roughly  $n/k$ . We show below that this is not always possible in high enough dimensions.

#### 5.1.1. A counterexample in high dimension

**Lemma 5.1.** *For any integer  $n > 0$ , there exists a set  $P$  of  $n$  points in  $\mathbb{R}^n$ , such that for any  $k < n$ , a  $k$ -density clustering of  $P$  must use at least  $n - k + 1$  centers.*

*Proof:* Let  $n$  be a parameter. For  $i = 1, \dots, n$ , let  $\ell_i = \sqrt{1 - 2^{-i-1}}$ , and let  $p_i$  be a point of  $\mathbb{R}^n$  that has zero in all coordinates except the  $i$ th one, where its value is  $\ell_i$ . Let  $P = \{p_1, \dots, p_n\} \subseteq \mathbb{R}^n$ . Now,  $d_{i,j} = \|p_i - p_j\| = \sqrt{\ell_i^2 + \ell_j^2} = \sqrt{2 - 2^{-i-1} - 2^{-j-1}}$ . For  $i < j$  and  $i' < j'$ , we have

$$\begin{aligned} d_{i,j} < d_{i',j'} &\iff 2 - 2^{-i-1} - 2^{-j-1} < 2 - 2^{-i'-1} - 2^{-j'-1} \iff 2^{-i'} + 2^{-j'} < 2^{-i} + 2^{-j} \\ &\iff \begin{cases} i = i' \text{ and } j < j', \text{ or} \\ i < i'. \end{cases} \end{aligned}$$

That is, the distance of the  $i$ th point to all the following points  $S_{i+1} = \{p_{i+1}, \dots, p_n\}$  is smaller than the distance between any pair of points of  $S_{i+1}$ .

Now, consider any set of centers  $C \subseteq P$ , and let  $c = p_i$  be the point with the lowest index that belongs to  $C$ . Clearly,  $P_c = (P \setminus C) \cup \{c\}$ ; that is, all the non-center points of  $P$ , get assigned by the clustering to  $c$ , implying what we want to show. ■

## 5.2. Algorithms

### 5.2.1. Density clustering via nets

**Lemma 5.2.** *For any set of  $n$  points  $P$  in  $\mathbb{R}^d$ , and a parameter  $k < n$ , there exists a  $k$ -density clustering with  $O(\frac{n}{k} \log \frac{n}{k})$  centers (the  $O$  notation hides constants that depend on  $d$ ).*

*Proof:* Consider the hypercube  $[-1, 1]^d$ . Cover its outer faces (which are  $(d-1)$ -dimensional hypercubes) by a grid of side length  $1/3\sqrt{d}$ . Consider a cell  $C$  in this grid — it has diameter  $\leq 1/3$ , and it is easy to verify that the cone  $\phi_C = \{tp \mid p \in C, t \geq 0\}$  formed by the origin and  $C$  has angular diameter  $< \pi/3$ . This results in a set  $\mathcal{C}$  of  $N = O(d^d)$  cones covering  $\mathbb{R}^d$ .

Furthermore, the cone  $\phi_C$  is formed by the intersection of  $2(d-1)$  halfspaces. As such, the range space consisting of all translations of  $\phi_C$  has VC dimension at most  $d' = O(d^2 \log d)$  [Har11, Theorem 5.22]. Let *slice* be the set that is formed by the intersection of such a cone, with a ball centered at its apex. The range space of such slices has VC dimension  $d'' = O(d + 2 + d') = O(d^2 \log d)$ , since the VC dimension of balls in  $\mathbb{R}^d$  is  $d + 2$ , and one can combine the range spaces as done above, see the book [Har11] for background on this.

Now, for  $\varepsilon = (k/N)/n$ , consider an  $\varepsilon$ -net  $R$  of the point set  $P$  for slices. The size of such a net is  $|R| = O((d''/\varepsilon) \log \varepsilon^{-1}) = d^{O(d)}(n/k) \log(n/k)$ .

Consider a point  $p \in P$  that is in  $R$ , and consider the cones of  $\mathcal{C}$  translated so that their apex is  $p$ . For a cone  $\phi \in \mathcal{C}$ , let  $\nu_\phi$  be the nearest point to  $p$  in the set  $\{R \setminus \{p\}\} \cap \phi$ . The key observation is that any point in  $P \cap \phi$  that is farther away from  $p$  than  $\nu_\phi$ , is closer to  $\nu_\phi$  than to  $p$ ; that is, only points closer to  $p$  than  $\nu_\phi$  might be assigned to  $p$  in the Voronoi clustering. Since  $R$  is an  $\varepsilon$ -net for slices, the slice ball( $p, \|\nu_\phi - p\|\)  $\cap \phi$  contains at most  $\varepsilon n = k/N$  points of  $P$ . It follows that at most  $k/N$  points of  $P \cap \phi$  are assigned to the cluster associated with  $p$ . By summing over all  $N$  cones, at most  $(k/N)N = k$  points are assigned to  $p$ , as desired. ■$

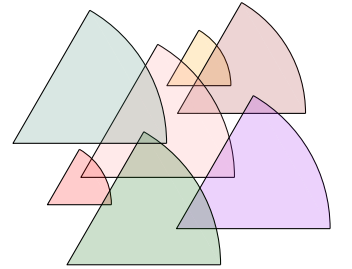
### 5.2.2. The planar case

**Lemma 5.3.** *For any set of  $n$  points  $P$  in  $\mathbb{R}^2$ , and a parameter  $k$  with  $1 \leq k \leq n$ , there exists a  $k$ -density clustering with  $O(n/k)$  centers.*

*Proof:* Consider the plane, and a cone  $\phi \in \mathcal{C}$  as used in the proof of Lemma 5.2. For a point  $p \in \mathbb{R}^2$ , and a radius  $r$ , let  $\phi_p$  denote the translated cone having  $p$  as an apex, and let  $s(p, r) = \phi_p \cap \text{ball}(p, r)$  be a slice induced by this cone. Consider the set of all possible slices in the plane:

$$S = \{s(p, r) \mid p \in \mathbb{R}^2, r \geq 0\}.$$

It is easy to verify that the family  $S$  behaves like a system of *pseudo-disks*; that is, the boundary of a pair of such regions intersects in at most two points, see figure on the right. As such, the range space having  $P$  as the ground set, and  $S$  as the set of possible ranges, has an  $\varepsilon$ -net of size  $O(1/\varepsilon)$  [MSW90]. Let  $\varepsilon = (k/N)/n$ , as in the proof of Lemma 5.2, where  $N = |\mathcal{C}|$ . Computing such an  $\varepsilon$ -net, for every cone of  $\mathcal{C}$ , and taking their union results in the desired set of cluster centers. ■





# Acknowledgments

N.K. would like to thank Anil Gannepalli for telling him about Atomic Force Microscopy.

## References

- [AS10] L.-E. Andersson and N. F. Stewart. *Introduction to the Mathematics of Subdivision Surfaces*. SIAM, 2010.
- [BGO07] J.-D. Boissonnat, **L. J. Guibas**, and S. Oudot. Learning smooth shapes by probing. *Comput. Geom. Theory Appl.*, 37(1):38–58, 2007.
- [Bli82] J. F. Blinn. A generalization of algebraic surface drawing. *ACM Trans. Graphics*, 1:235–256, 1982.
- [BQG86] G. Binnig, C. F. Quate, and Ch. Gerber. Atomic force microscope. *Phys. Rev. Lett.*, 56:930–933, Mar 1986.
- [Cla10] **K. L. Clarkson**. Coresets, sparse greedy approximation, and the frank-wolfe algorithm. *ACM Trans. Algo.*, 6(4), 2010.
- [CY87] R. Cole and C. K. Yap. Shape from probing. *J. Algorithms*, 8(1):19–38, 1987.
- [FG88] T. Feder and D. H. Greene. Optimal algorithms for approximate clustering. In *Proc. 20th Annu. ACM Sympos. Theory Comput.* (STOC), pages 434–444, 1988.
- [GIV01] A. Goel, **P. Indyk**, and **K. R. Varadarajan**. Reductions among high dimensional proximity problems. In *Proc. 12th ACM-SIAM Sympos. Discrete Algs.* (SODA), pages 769–778, 2001.
- [Gon85] T. Gonzalez. Clustering to minimize the maximum intercluster distance. *Theoret. Comput. Sci.*, 38:293–306, 1985.
- [Har11] **S. Har-Peled**. *Geometric Approximation Algorithms*, volume 173 of *Mathematical Surveys and Monographs*. Amer. Math. Soc., 2011.
- [HIM12] **S. Har-Peled**, **P. Indyk**, and R. Motwani. Approximate nearest neighbors: Towards removing the curse of dimensionality. *Theory Comput.*, 8:321–350, 2012. Special issue in honor of Rajeev Motwani.
- [HM06] **S. Har-Peled** and M. Mendel. Fast construction of nets in low dimensional metrics, and their applications. *SIAM J. Comput.*, 35(5):1148–1184, 2006.
- [Ind04] **P. Indyk**. Nearest neighbors in high-dimensional spaces. In J. E. Goodman and **J. O’Rourke**, editors, *Handbook of Discrete and Computational Geometry*, chapter 39, pages 877–892. CRC Press LLC, 2nd edition, 2004.
- [Kal12] B. Kalantari. A characterization theorem and an algorithm for A convex hull problem. *CoRR*, abs/1204.1873, 2012.
- [Man83] B. B. Mandelbrot. *The fractal geometry of nature*. Macmillan, 1983.

- [MB84] J. M. Mulvey and M. P. Beck. Solving capacitated clustering problems. *Euro. J. Oper. Res.*, 18:339–348, 1984.
- [MSW90] J. Matoušek, R. Seidel, and E. Welzl. How to net a lot with little: Small  $\varepsilon$ -nets for disks and halfspaces. In *Proc. 6th Annu. Sympos. Comput. Geom. (SoCG)*, pages 16–22, 1990.
- [PASG13] F. Panahi, A. Adler, A. F. van der Stappen, and K. Goldberg. An efficient proximity probing algorithm for metrology. In *Proc. IEEE Int. Conf. Autom. Sci. Engin. (CASE)*, pages 342–349, 2013.
- [SKG<sup>+</sup>09] R. M. Smelik, K. J. De Kraker, S. A. Groenewegen, T. Tutenel, and R. Bidarra. A survey of procedural methods for terrain modelling. In *Proc. of the CASA Work. 3D Adv. Media Gaming Simul.*, 2009.
- [Ski89] S. S. Skiena. Problems in geometric probing. *Algorithmica*, 4:599–605, 1989.
- [Ski97] S. S. Skiena. Geometric reconstruction problems. In J. E. Goodman and J. O’Rourke, editors, *Handbook of Discrete and Computational Geometry*, chapter 26, pages 481–490. CRC Press LLC, Boca Raton, FL, 1997.
- [Wik14] Wikipedia. Atomic force microscopy — wikipedia, the free encyclopedia, 2014.